LAWRENCE
LIVERMORE
NATIONAL
LABORATORY

# Analysis of Hypothetical Promoter Domains of DKFZp564A1164, NPHS1 and HSPOX1 Genes

S. S. Hammond

December 1, 2003

## Disclaimer

ANALYSISOFHYPOTHET ICALPROMOTERDOMAIN SOF

*DKFZP564A1164*, *NPHS1*AND *HSPOX1*GENES

UCRL-TH-201207

---

AUniversityThesisPresentedtotheFaculty

of

CaliforniaStateUniversity,Hayward

---

InPartialFulfillment

oftheRequirementsfortheDegree

MasterofScienceinBiologicalScience

---

By

ShaHammond

March,2004

# ANALYSISOFHYPOTHET ICALPROMOTERDOMAIN SOF

## *DKFZP564A1164*, *NPHS1* AND *HSPOX1* GENES

By

ShaHammond

Approved:                                        Date:

_____                    _____

_____                    _____

_____                    _____

_____                    _____

# ACKNOWLEDGMENTS

# TABLEOFCONTENTS

# FIGURESANDTABLES

## INTRODUCTION

Comparing the human genome to that of a related species, such as mouse, provides a unique perspective for identifying similarities and for finding the genes each sequence may encode. This approach has become a powerful method to identify sequences of specific function, such as gene regulatory activity (Loots *et al.*, 2000). Genome comparison works because the biologically essential features of a genome, such as genes and regulatory elements, are conserved through evolutionary pressure, while the non-essential elements readily acquire mutations and diverge between species. Deleterious mutations that occur within essential DNA are not conserved because they decrease the survival rate of the organism, while advantageous mutations, those that increase or preserve the survival rate of the organism, are conserved. This essential DNA is comprised of the protein incoding exons of genes and the regulatory sequences that control their activity (Hardison *et al*, 1997; Hood *et al*, 1993). The use of comparative sequence alignments is, therefore, an effective tool for providing confirmatory evidence of hypothetical genes by identifying candidate exons and regulatory elements, which can be difficult to ascertain through other predictive methods.

The comparative sequence analysis of human chromosome 19 (HSA19) and related regions in mouse highlighted the positions of more than 1300 genes and associated putative regulatory elements including promoters and enhancers (Dehal *et al.*, 2001). These elements are especially interesting because so little is known about them: for instance only 1871 promoters have been characterized out of the 30,000 total human genes (from the Eukaryotic Promoter Database http://www.epd.isb-sib.ch) (Praz *et al*,

2002).Inordertoconfirmornegatethefunctionalrelevanceofthislargenumberof

predictedregulatoryelements,wesetout    todevelopahighthroughputpipelinetotestfor

promoterandenhancerfunctioninculturedmammaliancells.

**SummaryandSignificanceoftheProposedResearch**

Thefocusofthismaster'sthesisprojectwastodevelopthebasicmethodsthat

willunderli eahighthroughputpipeline,andtousethesemethodstoinvestigatepotential

promoterelementsinaspecificgene    -richregioncontaininglociassociatedwithseveral

humandiseaseloci.Theregionoffocuswasa67kbsegmentofhumanchromosome

19q13.1(segmentofGenomicContig,GenbankaccessionnumberNT_011196.11),

containingthreegenes  *DKFZp564A1164* (*NLG1*), *NPHS1*and  *HSPOX1*(alsoreferredto

as *PRODH2*),figure1.HSA19waschosenasithasbeenthefocusofmyworkat

LawrenceLivermoreNational  Laboratories(LLNL)andthereisawealthofsequence

andexperimentaldataavailableforanalysisofthisverygenerichchromosome.

**Figure 1:A67kbsegmentofhumanchromosome19q13.1(smallportionofgenomiccontig,Genbank accessionno.NT_011196.11),containingthreegenesHSPOX1,NPHS1andDKFZp564A1164.**

## Genesofinterest

*NPHS1,HSPOX1* and*DKFZp564A1164* werechosenbecauseoftheirsimilar

expressionpatterns,andinadditiontheyarepartofalargerwell       -characterized gene-rich

regiononHSA19q13.1(LocusLink,NT_011196.11).Nephrin,the         *NPHS1*geneproduct,

isa1241 -residueputativetransmembranekidneyproteinoftheimmunoglobulinfamily

ofcelladhesionmolecules   (Kestila*etal*,1998)   .Thedisease ,congenitalnephritic

syndromeoftheFinnishtype,iscausedbymutationinthe         *NPHS1*gene,andexists

predominatelyinFinland  (Kestila*etal*,1998;Lenkkeri    *etal.* ,1999) .Itischaracterized

bymassiveproteinuria,detectableinute    robyalargeplacentaandmarkededema

(Hallman*etal*,1956)    .The *NPHS1*genehas29exonsandspans25.9kbinlength (GenbankaccessionNo.,NM_004646).

HSPOX1,alsoknownas    *PRODH2*,kidneyandliverprolinedehydrogenase (oxidase)2is  locateddownstreamof   *NPHS1*andhasaverysimilarexpressionpatternto thatofthe *NPHS1*gene. Theproteinencodedby*HSPOX1*    issimilarto *PRODH,*proline dehydrogenase(oxidase)1,amitochondrialenzyme,whichcatalyzesthefirststepin prolinecatabo lism.Thereissomeindicationthatheterozygousdeficiencyof      *PRODH*on HSA22maybeacauseofisolatedhyperprolinemia       (Goodman*etal.* ,2000) and schizophreniasusceptibility (Chakravarti,2002) . Theknown *HSPOX1* genes equence contains11exonsandisover13kbinlength(Genbankaccessionno.NP_067055). However,t hefunctionoftheproteinencodedby     *HSPOX1*  hasnotbeendetermined.

*DKFZp564A1164*isahypotheticalprotein(Genbankaccessionno.XP_048303) representedby acDNAisolatedfromhumanfetalbraintissue(AL136654)       (Wiemann*et al.*,2001) andretinoblastomacells(Genbankaccessionno.BC007312).Asrecentlyas January2003,Ihalmo *etal.* havedescribed *DKFZp564A1164* asanovelnephrin  -like gene( *NLG1*)encodingfiltrin,aproteinwithsubstantialhomologytohumannephrin.      The known *DKFZp564A1164*codingsequencecontains15exonsandis10kbinlength.       In additiontothefull   -lengthform,twoalternativelysplicedmRNAvariantswere discovered (Ihalmo*etal.* ,2003) . *NPHS1*and *DKFZp564A1164* aretranscribedin oppositedirectionsandthedistancebetweenthetranscriptionstartingpointsis approximately5 -kb,suggestingthatthesetwogenesshareacommonpromoterregion andenha ncers.

The mouse *Nphs1* gene promoter region has been previously reported and compared to human DNA by sequence alignment (Moeller *et al*, 2000). The corresponding *NPHS1* gene promoter region in human is conserved in sequence, as highlighted by our percent identity plot (PIP) in figure 2 and VISTA ([http://www-gsd.lbl.gov/vista/](http://www-gsd.lbl.gov/vista/)) alignment in figure 3. However, the precise locations of regulatory elements and start site for transcription of Nphs 1 have not been defined.

**Sequence comparison tools**

The percent identity plot (PIP) is one of the displays available from PipMaker ([http://bio.cse.psu.edu](http://bio.cse.psu.edu)), a site for comparing two long DNA sequences to identify conserved segments between species (Schwartz *et al*, 2000). A PIP shows the position in one sequence of each aligning gap-free segment and plots the degree of similarity between both species as dots or lines (similar to dotplot). For example, PIPMaker can align completed human sequence with homologous mouse DNA even if it is draft sequence, and reveal candidate regulatory elements as highly conserved regions that do not correspond to exons or predicted exons. Positions along the horizontal axis can be labeled with known features such as exons, repetitive elements and CpG islands (Figure 2).

**Figure 2:PIPcomparingaregioninhuman19q13.1andinmouse,highlightinghypothetical promoters(purple)andfirstexons (pink)predictedbytheFirstEFprogram(Davulurietal.,2001). Numbers1and2designateFirstEFpromoterpredictionsfor** *NPHS1***.**

VISTAisaprogramforvisualizinggglobalDNAsequencealignmentsofarbitrary

length.Itwasdesignedtovisualizelongse quencealignmentsofDNAfromtwoormore

species,suchashumanandmouse,withannotationinformation (Bray,2003;Dubchak *et*

*al.*,2000;Mayor *etal.* ,2000) .VISTAiseasilyconfigurable,allowingthevisualization

ofalignmentsofvari ouslengthsatdifferentlevelsofresolution.Infigure3thex -axis

representsbasesequencesandthey -axisrepresentspercentidentityofconserved

sequencesintheformofgraphicalpeaks.AsonecanseesomesegmentsofDNAare

highlyconservedwhe reasotherregionsareverydissimilarbetweenthehumanand

mouse. Different sequence features such as exons and UTR's are denoted by color coding (Figure 3).



**Figure 3: VISTA comparing a region in human 19q13.1 and consensus sequence in mouse, peaks represent conserved sequence.**

The human *NPHS1* promoter has not been characterized in laboratory experiments. The hypothetical promoter for this region is 4 to 5 kb upstream of the currently known first exon of the *NPHS1* human gene, suggesting that there is another undiscovered upstream exon for this gene. In fact, many known gene sequences are not

complete,inthesensethattheyincludethefullprotein -codingsequencebutdonot

containacompletesetofnon -coding5'exons (Davuluri*etal*,2001) .Inaddition,a

growingbodyofdatasuggeststhatmanygenesusealternativestartsitesandpromoters

indifferenttissues (Asnagli*etal.* ,2002) .Identifyingstartsitesandallpromotersusedby

HSA19genesisthegoalofalargerstudyintheStubbslaboratory,andthismastersthesis

wasdesignedasafocusedpilotstudytotestmethodsandapplythemtoanalysisofthe

*NPHS1*generegion.

**FirstEFannotation**

Thepositionsofpromoters(purple)and firstexons(pink),whichwerepainted

ontofigure2,werepredictedbyFirst –exonfinder(FirstEF,

http://www.cshl.org/mzhanglab),aprogramdevelopedbyM.Zhangandcolleaguesat

ColdSpringHarborLaborat ory (Davuluri*etal.* ,2001) .Weareworkingincollaboration

withM.ZhangandZhenyuXuan(ColdSpringHarborLaboratory)toconfirmthe

FirstEFpredictionsinourlaboratoryusingexperimentalmethods .

FirstEFconsistsofasetofdiscri minantfunctionsdesignedtofindpotentialfirst

splice-donorsitesandCpG -islandrelatedandnon -CpG-islandrelatedpromoterregions.

FirstEFdecideswhethertheintermediateregioncouldbeapotentialfirstexonand

upstreampromoterbasedonthisset ofquadraticdiscriminantfunctions.Forexamplethe

regionslabeled1and2infigure2arepredictedbyFirstEFtobepromotersforthe

*NPHS1*gene,andregion1isalsopredictedtobeapromoterfor *DKFZp564A1164*

(althoughinthereverseorientation) . Noindependentpromoterwaspredictedfor

*HSPOX1* byFirstEF ,however,thesimilarityinexpressionpatternsbetweentheNPHS1    ,

*DKFZp564A1164*and  *HSPOX1*genesledustohypothesizethatthesegenesmaysharea

singlepromoter.  Thepotentialsharingofas   inglebi -directionalpromoterregionbythree

neighboringgenesmadethisregionespeciallyinterestingtoanalyze.

**Bioluminescentreporterassay**

Toanalyzethefunctionalityofthispossiblebi      -directionalpromoteraswellas

otherputativepromotersinthe      *NPHS1*generegion,atransientluciferasereporterassay

wasused.Bioluminescentreporterassayshaveawiderangeofapplicationsincludingthe

functionalanalysisofpromotersandenhancers,andithasbeendemonstratedthatthese

systemsprovide  reliablereproducibleresults  (Parsons,2000;Sherf,1996)   .

TheDual -luciferasereportersystem(PromegaCorporation)utilizesfireflyand

Renillaluciferaseinaco    -reportersystemwhereRenillaisaninternalcontrolallowingfor

normalizationofthefireflyluciferasedata.Inthisstudy,theregionspredictedtobe

promotersbyFirstEFwereplacedintovectorsthatexpressfireflyluciferasewhen

borderedbyafunctioningpromoterandtransfectedintotheappropriateeukaryoticcell

lines.

Preliminaryexpressiondatawereusedasaguideinchoosingtheappropriatecell

linesforourtransientreporterassaystudies.Expressionprofilesforthesegeneswere

obtainedfromanumberofsourcesincludingtheGenbank'sSAGEandESTdatabases

(serialanalysisofgeneexpression,andexpressedsequencetag,respectively,

http://www.ncbi.nlm.nih.gov/sage),adatabaseofgeneexpressionusingmicroarrays

calledGeneExpressionAtlas( [http://expression.gnf.org/cgi-bin/index.cgi](http://expression.gnf.org/cgi-bin/index.cgi)),andtissue

section *insitu* hybridizationanalysisthatwasperformedatLawrenceLivermoreNational

Laboratory(LLNL).Finally,themostlikelycandidatecelllinesweretestedfor

expressionofthegenesofinte   restusingRT -PCR(reversetranscription -polymerasechain

reaction)andgenespecificprimers.

## ThesisObjective

Thisthesis'primaryobjectivewastousecomparativesequenceanalysisprograms

suchasPipMakerandVISTA,inadditiontothecomputational         programFirstEF,to

identifypotentialpromotersandenhancersforthreegenesinthe      *NPHS1*region,andto

testtheseregulatoryelementsinculturedmammaliancelllinesusingtransiently

expressedluciferasereporterconstructs.Additionally,determinin    gthefirstexonsfor

*NPHS1*,*HSPOX1* and  *DKFZp564A1164,*includingpotentialalternativestartsiteslinked

todifferentpromoterswasattemptedandresultssequenced.Overalltheaimhasbeento

testthehypothesisthatasinglebi    -directionalpromoterwas beingsharedby *NPHS1,*

*DKFZp564A1164*and  *HSPOX1,*threeneighboringgeneswithsimilarexpression

patterns,andtoestablishthetechnologyandmethodsforahighthroughputassayof

promoterandenhancerelements.

## MATERIALSANDMETHODS

**SequenceComparis ons**

An845kbcontigfromhumanchromosome19(Genbankaccessionno.
NT_011296)andrelatedregionsinmouse(Genbankaccessionnos.AC087141and
AC020839)werecomparedusingthePipMakerprogram
(http://bio.cse.psu.edu/pipmaker) (Schwartz*etal*,2000) .Inspeciesthatdiverged100 -
300millionyearsago,suchashumanandmice,exonsandgeneregulatoryelementsare
detectableassimilarsequences.Thesecanbevisualizedonapercentidentityplot(PIP),
whichshowsthepositioninon esequenceanddegreeofsimilaritybetweenthealigning
sequences (Schwartz*etal*,2000) .IncollaborationwithM.Zhang(ColdSpringHarbor
Laboratory),FirstEFpredictionswereusedtoanalyzethesequence,andregions
predictedtobe hypotheticalpromotersbyFirstEFwerefurtheranalyzedforpromoter
activity.

**CellCulture**

HumanandmousecelllinesfromAmericanTypeCultureCollection(ATCC)
wereculturedinmediaandserarecommendedbyATCCandcontaining100I.U./mlof
penicillin,100µg/mlstreptomycinand 2mMofL -glutamine.Growingcultureswere
housedinacellcultureincubatorat37 °Cwith5%CO$_2$ orasrecommended.We
preliminarilyselectedthecelllinesbasedonpubliclyavailableSAGEexpressiondata
(NCBI)forHSA19genes,forgrowthcharacteristi cs,fortransfectionassayperformance

(basedonourownresultsandpublisheddata),andtorepresentawidevarietyofcell

typesandtissues.

**AnalysisofcDNA**

Todeterminewhichcelllinesexpressthegenesofinterest,RNAwascollected

fromthemost likelycellcandidatesbasedonexpressiondataobtainedonpublic

databasesorpreviousstudies,andcDNAwasproducedviaRT-PCRusingthe

RNAqueouskit(AmbionInc.).CellsweregrownasrecommendedbyATCCuntilthey

reachedayieldof1 $\times 10^5$to10 $^8$,thenthecellswerecollectedandstoredinRNAlater

(AmbionInc.)untilcDNAwasmade.Primersweredevelopedthatspecificallyamplified

the3'endsofthecDNAofinterest,andstandardPCRwasperformedusingPerkinElmer

reagentsonanMJResearchthermocycler.Primersequencesarelistedintable1Ainthe

Appendix.Ifabandwasproducedoftheexpectedsize,thenthatcelllinewasconsidered

toexpressthegeneandwasusedinsubsequenttransfectionassayexperiments.

**5′EndTranscriptVerification**

Inthecaseof *HSPOX1*whereFirstEFandothermethods,suchasthepresenceof

CpGislandsorGATAandTATAboxes,didnotpredictapromoterandfirstexon,5′

SMARTRACE(BDBiosciencesClontech)wasperformedtoverifythepositionofthe

firstexon.SMARTRACEincorporatesa switchingm echanism atthe5′endofan RNA

transcriptcoupledwithRACE(rapidamplificationofcDNAends)toisolatethecomplete

5′endsequenceofatargetgene.Additionally5′SMARTRACEwasperformedon

*NPHS1* as FirstEFpredicted 2 first exons for this gene. Often it is the case that the transcription start site is upstream from the start ATG codon in an untranslated initial exon. It was the hoped that 5′ RACE would help to identify any possible untranslated initial exons, and also to establish the sequence of the proximal promoter. After performing 5′ SMART RACE the PCR product was subcloned into a TA vector (Invitrogen Corp,) and sequenced using vector primers [m13(-20) and m13] on an ABI Prism 377 sequencer.

## ConstructDevelopment

### Vectorpreparation

The pGL3 enhancer or promoter vectors (Promega Corporation) were double digested overnight with the appropriate restriction enzymes (MluI and BglII or KpnI and BglII from New England Biolabs, Inc.) for directional subcloning, then the vector was dephosphorylated to prevent recircularization using alkaline phosphatase from calf intestine (New England Biolabs, Inc.). Following which the vector was purified from an agarose gel using a Qiagen kit and eluted in TE. At est of the vector's re-ligation efficiency was performed by transforming Electromax cells (Gibco Invitrogen Corporation) and growing on an LB/AMP plate overnight. Vectors were considered good if less than 75 colonies grew.

### Insertpreparation

Primers were designed that flanked the hypothetical promoters and contain restriction sites at the 5′ end complementary to the sites in the vector's multi-cloning site.

ThenPCRwasperformedandasmallaliquotrunonageltodeterminethatthePCR

worked.ThePCR   productwastreatedwithKlenowfragment(NewEnglandBiolabs,

Inc.)tofillin3   ′recessedends,andthenthePCRproductwasdoubledigestedwiththe

appropriaterestrictionenzymesandgelpurified.


Ligation

ThepGL3 -Enhanceror  -Basicvectorandinse   rtwereligatedwithT4DNAligase

(NewEnglandBiolabs,Inc.)overnightusinganexcessofinsert.Electromaxcellswere

transformedwiththeligationproductandplatedovernightonLB/AMPafteroutgrowth

for1hourinLBonly.Colonieswerescreenedv        iaPCRusingvectorspecificprimers,

andthosethatcontainedtheinsertweregrowninLB/AMPovernightandisolatedusing

theQiagenHighSpeedMidiprep.Analiquotoftheisolatedconstructswasconfirmedby

restrictiondigestionornestedPCRandlate      rsequenced.


**TransfectionAssays**

DualLuciferaseTransfectionAssays(PromegaCorporation)wereperformedto

determineifthepredictedpromotersfunctioned      *invitro* .Bioluminescentreporterassays

havebeendemonstratedtoprovidereliablereproducible       resultsforthefunctional

analysisofpromotersandenhancers      (Parsons,2000;Sherf,1996)   .Promoterassayswere

performedusingthepGL3   -Enhancervectorandinternalcontrolco      -reporter,pRL -TK

(PromegaCorporation).Promoterandenha       ncerassayswereperformedusingthepGL3    -

Basicvectorandthesameinternalcontrolco        -reporter.

<u>pGL3-EnhancerVector</u>

ThepGL3 -Enhancervectorcontains *luc*+cDNA,whichencodesmodifiedfirefly luciferase,amultiplecloningregionupstreamof *luc*+for insertionofthepromoter element,andanSV40enhancerlocateddownstreamof *luc*+ .TheSV40enhanceraidsin theverificationoffunctionalpromoterelementsbyincreasingthelevelsof *luc*+ transcription.

<u>pGL3-BasicVector</u>

ThepGL3 -Basicvectorconta ins *luc*+cDNA,whichencodesmodifiedfirefly luciferase,andamultiplecloningregionupstreamof *luc*+forinsertionofthe promoter+enhancerelement.ThepGL3 -BasicvectordoesnotcontainanSV40enhancer orpromoterinordertodeterminethepresence ofafunctionalenhancerandpromoterin theexperimentalconstruct.

<u>pRL-TKVector</u>

ThepRL -TKvectorisaninternalcontrolreporterintendedtobeusedin combinationwithanyexperimentalreportervectortoco -transfectmammaliancells.The pRLreporte rvectorcontainsacDNA( *Rluc*)encodingRenillaluciferase,whichwas originallyclonedfromthemarineorganism *Renillareniformis* (seapansy).ThepRL -TK vectoralsocontainstheherpessimplexvirusthymidinekinase(HSV -TK)promoterto

providelowto moderatelevelsofRenillaluciferaseexpressioninco      -transfected

mammaliancells.

<u>pGL-ControlVector</u>

ThepGL -ControlvectorcontainstheSV40promoterandenhancersequences,

resultinginstrongexpressionof      *luc*+inmanymammaliancelltypes.Thisi      susefulin

monitoringtransfectionefficiencyingeneralandisaconvenientinternalstandardfor

promoterandenhanceractivity.ThespecifictranscriptionalactivityofpGLvectors

variesfordifferentcelltypesandthepGL      -Controlvectorcanhelpde   termineactivityto

beexpectedbyastrongpromoter.

<u>Transfection</u>

HumancelllinesHepG2and293,determinedtoexpressthegeneofinterestby

analysisofcellularcDNAwithgenespecificprimers,wereplatedina96wellformat.

Onehundredmicroliter sofcellswereplatedinOpti      -MEM(GibcoBRL)at1      $\times 10^4$ cells

perwellinthecenter60wells.Theouterwellswerefilledwith100µlofPBStoprevent

drying.Twenty -fourhourslaterthecellsweretransfectedwiththevectorsvialipofection

accordingtotheFugene6TransfectionReagentprotocol(RocheMolecular

Biochemicals).Foreachwell,5.52      $\times 10^{-14}$ molesofexperimentalvector(pGL3plus

insert)and50ngofco -reporter(pRL)weremixedwith0.3      –1.8µlofFugenein5µl

Opti-MEMandaddedtoeac   hwell.Moleswerechosenasthemeasuringunitforthe

experimentalconstructstohelpensureanequalamountofeachconstructwasdelivered

Constructs ranged in size from 5.5 to 8.0 Kb. The plates were then incubated for approximately 24 hours before   the Dual Luciferase assay was performed.

DualLuciferaseAssay

The assays were conducted according to Promega's Dual Luciferase $^{®}$Reporter 1000 Assay system. The media from each well was removed, and 20µl of PLB lysate was added to each well. The plate w      as then incubated at 37 ˚C with shaking until the cells were completely lysed. Then 100µl of LARII was added, activating any firefly luciferase generated by the pGL3 construct. The RLU (relative light units) of firefly luciferase was measured using a Packar   d LumiCount microplate luminometer set to a 5 second read time and a 1 second delay between reads. Then 100µl of Stop&Glo reagent was added to each well and the RLU of   *Renilla* luciferase from the pRL   -TK co -reporter was measured.

DataAnalysis

Each constr uct was transfected in 3 wells and each well was measured in triplicate. Firefly measurements were averaged for each construct, and Renilla measurements were also averaged per construct in the same manner. The negative control treatment contained pGL   -Enhancer or pGL3 -Basic construct without insert co   -transfected with pRL   -TK. The positive control treatment contained pGL      -Control and was used as a comparison to maximum expression. The measurements for these were averaged in the same way. The change in fol      d activity was determined by dividing the sample ratio by the negative control ratio (firefly avg. RLU divided by Renilla avg.

RLU).Constructsthatcausedanincreaseinfoldactivityabovethenegativecontrols
wereconsideredtocontainaworkingpromo ter.


**rVISTAAnalysis**

AnalysisusingrVISTA( http://www-gsd.lbl.gov/vista/)wasperformedonallthe
constructsdevelopedandtransfectedtobetterunderstandwhichtranscriptionfactor
bindingsitesmaybec ontainedwithintheconstructsandthereforewhichtranscription
factorsmaybeactingonthesequences.rVISTAisacomputationaltoolfor
comparativesequence -baseddiscoveryoffunctionaltranscriptionfactorbindingsites
(TFBS) (Loots,2002).

Morespecifically,rVISTAenablesthehighthroughputdetectionofcis -regulatory
elementsbycombiningclusteringandanalysisofconservedinterspeciessequenceto
maximizetheidentificationoffunctionalsites.InitiallyrVISTAalignshuman and
mousesequencesusingAVID,aglobalalignmentprogram.Thenpotentialtranscription
factorbindingsitesarepredictedbyMatch ™programbasedonTRANSFAC
Professionallibrary5.3.AfterfindingalltheTFBSineachspeciesindependently,the
siteswherecorepositionscorrespondinbothspeciesareselectedasalignedsites.
Finally,onlythealignedtranscriptionfactorbindingsitesthatarefoundwithinconserved
human-mousesequenceatalevelof80%ormoreareselectedbyrVISTAasprobable
transcriptionfactorbindingsites.

**RESULTS**

**PreliminaryExpressionData**

Expressiondatawasusedasaguideinchoosingtheappropriatecelllinesforour transientreporterassaystudies.Expressionprofilesforthesegeneswereobtainedfroma number ofsourcesincludingtheGenbank'sSAGEandESTdatabasesaswellasGene ExpressionAtlas'microarraydatabase.Additionally,tissuearrayanalysiswasalso performedatLLNL.

<u>SAGEandEST</u>

TheSAGEdatabaseusesatechnique,whichquantifiesa"tag"that representsthe transcriptionproductofagene.Thenumberoftimesaparticulartagisobserved providestheexpressionlevelofthecorrespondingtranscript.Thehistogramdenotes expressionlevel.UsingtheSAGEhistogramasaguide,thestrongestex pressionof *NPHS1*wasfoundinthekidney,brain,mammaryglandandtestistissues. SAGEexpressiondataalsoshowed *HSPOX1*tobeexpressedinkidneyandnormalliver tissue.

TheESTdatabaseshowed *NPHS1*expressionintheendometrium, adenocarcinomacell lineandIsletsofLangerhans;and *HSPOX1*expressioninliver, spleenandkidney.BothGenbank'sSAGEandESTexpressiondatashowed *DKFZp564A1164*tobeexpressedinbrain,germcells,kidneyandlung.

GeneExpressionAtlas

GeneExpressionAtlas( http://expression.gnf.org/cgi-bin/index.cgi)microarray databaseonlycontaineddatafor2ofthegeneswetested( *NPHS1, HSPOX1)*.For *NPHS1*strongpositives(3 ×abovemedian)werenotedintheDOHH2, lymphomaandB celllines,andinthekidney,pituitaryandpancreas(10 ×abovemedian)tissues.For *HSPOX1*positives,all10 ×abovemedian,werenotedin3tissues:kidney,fetalliverand liver.Therewerenopositives3 ×abovemedianfor *HSPOX1*. (Append ix,Figures1Aand 2A)

TissueArrayResults

Forthetissuearrayanalysis *NPHS1, HSPOX1*and *DKFZp564A1164*geneswere hybridizedtohumantissuearrayslidesbyX.LuandE.Wehri,atLLNL.T7mRNA probesweremadeusingthemRNAsequenceofeachgeneand orderedfromLife Technologies(GibcoBRL).Thesequenceofeachprobeislistedintable2Ainthe Appendix.TheprobeswerethenlabeledwithdigandhybridizedtoMaxArraynormal humantissueslides(ZymedLaboratories,Inc.)usingstandardprotocols.

Theresultsingraph1bellowindicatepositiveexpressioninseveraltissues, ranginginlevelfrom5to15(arbitraryvalues).Forexample,allthreegenesarehighly expressedinthetestisandovaryandmoderatelyexpressedinthekidneytubules. *DKFZp564A1164*and *HSPOX1*wereexpressedinliverwhileonly *DKFZp564A1164*was expressedinlung,heartandcolon.Noneofthegeneswereexpressedinthespleenor skeletalmuscle.Table1AintheAppendixshowsallthetissuestestedandtheexpression

result.Figures3A,4Aand5A(Appendix)arepicturesofthehybridizationresultson

livertissue.Fromtheseresultsonecanseethatintheliver *NPHS1*isveryweakly

expressed, *HSPOX1*ismoderatelypositiveand *DKFZp564A1164*ispositiveand

expressingth egeneinspecificcellsoftheliver.

Theresultsofthetissuearrayexperimentsareuniqueinthattheycanshowthe

typeofcellwithinatissuethatisexpressingthegene.Moreoftenthannotageneis

expressedinaspecificcelltypeinthetissue andnotthewholetissue.Inkidneyfor

example,theexpressionof*DKFZp564A1164* and *HSPOX1*wereonlyseeninthecells

liningthetubules(datanotshown).Forthisreasonthedataarenotalwaysthesameas

otherexpressionstudieswhereresultsfroma wholetissueorindividualcelllineare

examined.



**Graph1:Comparisonofpositivetissuehybridizationresults.**

**cDNA Analysis and RT-PCR Results**

PCR was performed on cDNA made from RT-PCR of the individual cell lines or RT-PCR of polyA⁺ RNA purchased from BD Biosciences. Commercial cDNA from Clontech was also tested as a control. Samples were run on a 1.2% agarose gel containing 25 µg of ethidium bromide at 90V for 30 minutes. Gels were then imaged using the AlphaImager 2000. Cell lines were considered to be expressing the gene if a band of the expected size was seen on an agarose gel.

All the primers were designed from the 3' end of the cDNA to span an intron so that a size difference could be visualized between genomic and cDNA. *NPHS1* cDNA size was 273bp and genomic DNA was 517bp, *DKFZp564A1164* cDNA size was 391bp and genomic DNA was 3Kb, likewise *HSPOX1* cDNA size was 306bp and genomic was 3Kb. As a positive control primers amplifying β-actin were used, and a PCR reaction lacking any template was used as the negative control.

*NPHS1* and *DKFZp564A1164* were found to be expressed in several human cell lines including 293, MDA MB-436, and PANC-1, as seen in figure 4 and table 1. Alternatively, expression of *HSPOX1* was only found in 2 of the human cell lines tested, Capan-1 and HepG2. Although *HSPOX1* is expressed in human kidney tissue, there was no indication of expression in the human kidney cell line 293 using this method **.**This may be due to the fact that expression data from RT-PCR of individual cell lines often differs from tissue analysis due to the difficulty of maintaining the tissues differentiated function *in vitro* (Mather & Roberts, 1998).

**Figure 4:PCRofcDNAfrom2 93,PANC -1,HepG2andMDA -MB-436celllines.Bandsrepresenting** *NPHS1* **and** *DKFZp564A1164* **expressionareseenin293,PANC -1andMDA -MB-436.** *HSPOX1* **and** *DKFZp564A1164* **expressionareseeninHepG2.Actinisapositivecontrol.**


BasedonthisinformationHepG 2and293celllineswerechosentobeusedinthe

transienttransfectionluciferaseassaysbecauseoftheirclearunambiguousresultsand

*HSPOX1* and *NPHS1* whichmaybesharingabi -directionalpromoter,aredifferentially

expressedinthesecelllines.LN CaP.FGC,whichdidnotshowexpressionofanyof

thesegeneswasusedtotestluciferaseassayresultsinanon -expressingcellline.

**Table1:RT -PCRExpressionData**

| CellLineor Tissue(human) | Tissue | Gene | cDNApresent |
|---|---|---|---|
| 293 | kidney | 1164** | yes |
| 293 | kidney | NPHS1 | yes |
| 293 | kidney | HSPOX1 | no |
| Capan-1 | pancreas | All* | yes |
| HelaS3 | cervix | HSPOX1 | no |
| HelaS3 | cervix | 1164 | faintband |
| HelaS3 | cervix | NPHS1 | yes |
| HepG2 | liver | 1164** | yes |
| HepG2 | liver | HSPOX1 | yes |
| HepG2 | liver | NPHS1 | no |
| IMR-32 | neuroblast | All* | no |
| Jurkart | leukemia,T -cell | All* | no |
| k562 | leukemia | All* | no |
| LNCaP.FGC | prostate | All* | no |
| MDA-MB-436 | breast | 1164** | yes |
| MDA-MB-436 | breast | HSPOX1 | no |
| MDA-MB-436 | breast | NPHS1 | yes |
| MDA-MB-453 | mammary | All* | no |
| PANC-1 | pancreas | 1164** | yes |
| PANC-1 | pancreas | HSPOX1 | no |
| PANC-1 | pancreas | NPHS1 | yes |
| commercialRNA | kidney | HSPOX1 | yes |
| commercialRNA | kidney | 1164** | no |
| commercialRNA | kidney | NPHS1 | yes |
| commercialcDNA | brain/testis | NPHS1 | yes |
| commercialcDNA | brain/testis | 1164** | yes |
| commercialcDNA | brain/testis | HSPOX1 | no |

*DKFZp564A1164,NPHS1,H SPOX1
**DKFZp564A1164

## 5′EndTranscriptVerification

FiveprimeRACE(SMARTRACE,BDBiosciencesClontech)wasperformedto

verifythepositionofthefirstexonforboth *HSPOX1*and *NPHS1*.Oftenitisthecase

thatthetranscriptionstartsiteisupstr eamfromthestartATGcodoninanuntranslated

25

initial exon. The 5′RACE experiment served to identify a possible untranslated initial exon, and therefore also to establish the position of the proximal promoter. After performing 5′RACE the PCR product was subcloned into a TA vector (Invitrogen Corp,) and sequenced using vector primers on an ABI Prism 377 sequencer.

As starting materials commercial liver and kidney polyA⁺ RNA from BD Biosciences Clontech were used. These RNAs were initially tested for the presence of the *HSPOX1* and *NPHS1* cDNA using the same primers designed to test the cell line RNA in the cDNA analysis method above.

The initial results from the 5′RACE were inconclusive. After several separate SMART RACE experiments, the 5′region of both *HSPOX1* and *NPHS1* have still not been identified. Not only were no new untranslated first exons identified, but also the currently accepted 5′end of these genes could not be verified using this method. The positive control provided with the kit was used in conjunction with these experiments and did produce the expected results.

In performing the 5′RACE experiments on the *HSPOX1* gene it was noted that the gene's first and second exons matched to multiple sites in the genome using NCBI BLAST, and when aligning the human and mouse mRNA sequences, it was found they do not form a consensus sequence alignment until base pair 298 in human which is equivalent to amino acid 77. Even when choosing unique primers from the consensus region, the 5′end of the gene was not found using SMART RACE. RACE products were generated but sequence did not correspond to any sequence from this genomic region.

Although the *NPHS1* gene is well characterized, and the mouse consensus region matches well, the 5′ end of the gene was not established using SMART RACE. The sequence of *NPHS1* RACE products matched ring finger/*DORFIN*, crystalin/ *CRYL1* glutathion/*GSTA2* and ribonuclease/ *PARN*, indicating that false priming was generating artifacts from abundant RNAs in the sample.

Primer design was of critical importance in these experiments. The 30 bp primers designed for these experiments had to match the gene of interest exclusively; if any part also matched a different area of the genome one risked amplifying both regions. Careful screening of not just the whole gene specific primer, but small segments of the primer was therefore necessary. BLAST searches revealed that the exons of these genes (*HSPOX1* and *NPHS1)* are littered with small sequence segments of 10 to 20 bp in length th at match other regions of the genome, making it difficult to find 30 bp gene specific primers for the SMART RACE experiments (primer sequences: Table 5A, Appendix). These repeat sequences most likely explain the failure of RACE to generate *NPHS1* and *HSPOX1* specific transcripts.

Under these circumstances the published 5' ends are probably the true ends of these transcripts, at least in the cell types tested. Since certain promoters may operate only in specific tissue types, it is possible that exhaustive R ACE in many tissues would have eventually yielded additional 5' sequences. However, such a search was beyond the scope of this study.

**TransfectionAssays**

DualLuciferaseTransfectionAssays(PromegaCorporation)wereperformedto determineiftheFirst  EF-predictedpromotersfunctionedaspromoters   *invitro* . Bioluminescentreporterassayshavebeendemonstratedtoprovidereliablereproducible resultsforthefunctionalanalysisofpromotersandenhancers       (Parsons,2000;Sherf, 1996).P romoterassayswereperformedusingthepGL3       -Enhancervectorandinternal controlco -reporter,pRL -TK(PromegaCorporation).Promoterandenhancerassayswere performedusingthepGL3   -Basicvectorandthesameinternalcontrolco       -reporter.

PreliminaryT ransfectionData: *XRCC1*

Inordertodeterminetheeffectiveness      ofthePromega'sDualLuciferaseAssay the *XRCC1*genewasshotgunsubclonedintothepGL       -Enhancervector.Byaligningthe baboon*XRCC1* promotersequence(Genbankaccessionno.      AF019114),whic hhadbeen previouslyclonedandcharacterizedbyZhou      *et.al,* withhuman(Genbankaccessionno. L34079)andmouse(Genbankaccessionno.      L34078)usingmVISTAwewereable visualizethehumanpromoterregion(Figure5)      (Zhou&Walter,1998   ).

Thenwebcutter(  http://www.firstmarket.com/cutter/cut2.html)wasusedto determinewhichrestrictionenzymewouldbebesttouse,andthehumanclone(Genbank accessionno. L34079)wasdigest edwithSacI.TheSacIdigestresultedinseven fragments,allofwhichwereshotgunsubclonedintothepGL3       -Enhancervector. Colonieswereisolatedthathad3.7kb,3.8kb,3.9kband7.9kbinserts,andthesewere testedusingPromega'sDualLuciferaseAssa     y.Promega'spGL3 -Control,whichcontains

aSV40constitutivepromoter,wasusedasapositivecontrol,andanemptypGL3         -

Enhancerwasusedasthenegativecontrol(Figure6).Thefragmentswerealsosequenced

andpositionalverified.

Theresultsshowedth   atonlythevectorcontainingthepromoterworked.The

othershotgunsubclonedsequenceshadvaluessimilartothenegativecontrol

demonstratingthattheassaydoesnottypicallygeneratefalsepositives.



**Figure 5:Alignmento fhuman,baboonandmousesequenceusingmVISTA(Genbankaccessionnos. L34079,AF019114,L34078respectively).**



**Figure 6:LuciferaseAssayof   *XRCC1*shotgunsubclonedfragments.**

ConstructDesign:  *NPHS1, HSPOX1* and *DKFZp564A1164*

ThepGL3 -Enhanceror -Basicvectors(PromegaCorporation)weredouble digestedwithrestrictionenzymesKpnI/BglIIorMluI/BglII(NewEnglandBiolabs, Inc.)fordirectionalsubcloningandligatedwithinsertsthatweredouble        digestedinthe samemanner.Therestrictionenzymeswerechosenbaseduponascreenofeachinsertto determinewhichrestrictionenzymesitestheydidnotcontain(Webcutter2.0,copyright 1997MaxHeiman,   [http://www.firstmarket.com/cutter/cut2.html](http://www.firstmarket.com/cutter/cut2.html)).Figure7isanexample ofinsertdesignandfigure8showstheregioneachvectorwasdesignedfromandnames eachweregiven.Table2givesadditionalinformationabouteachconstructincludin        g sizeandregionofcosmidR33502(Genbankaccessionno.AC002133)theywerecloned from.TheprimersforeachinsertweredesignedwithaBglII,KpnIorMluIsiteaddedto the5′endaccordingtorecommendationsinNewEnglandBiolabstechnicalliteratu        re. Primersequencesarealllistedintheappendixtable4A.

**Figure 7:PromoterandenhancerinsertsfordirectionalsubcloningintopGL3vectors.**



**Figure 8:Nameofeachinsertandregionitwasdev        elopedfromalongthepipplot.**

**Table2:Summaryofconstructs**

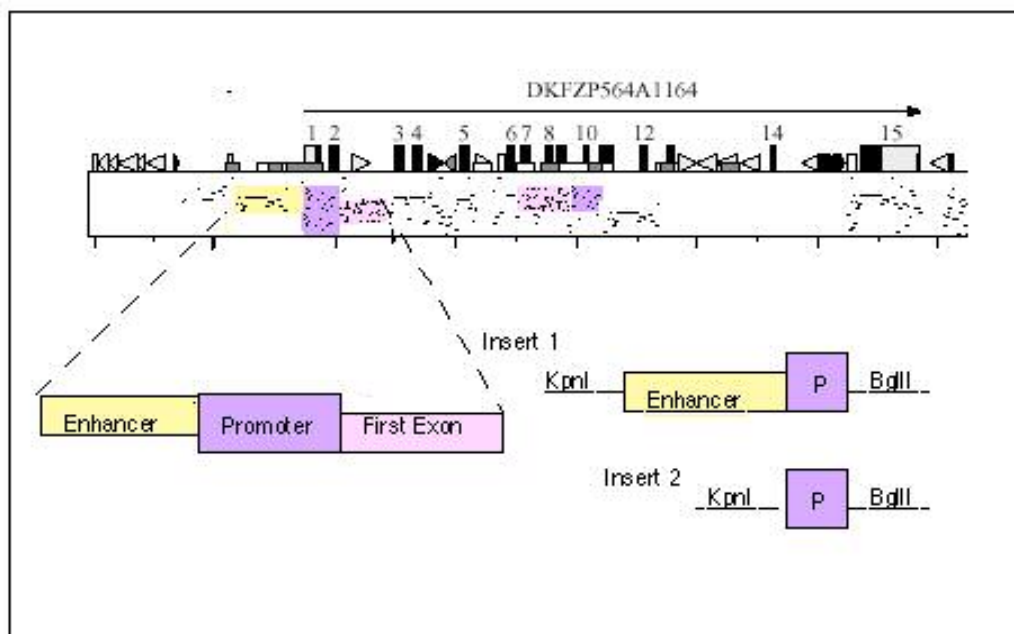| Name | Insert size | Type | Vector | Digest | Regionincosmid R33502 (AC002133) |
|---|---|---|---|---|---|
| n1-Basic | 1.2kb | promoter+enhancer | pGLBasic | BglII/KpnI | 37619-36392 |
| n2-Enhancer | 570bp | promoter | pGLEnhancer | BglII/KpnI | 36964-36392 |
| n3-Enhancer | 619bp | promoter | pGLEnhancer | BglII/MluI | 32729-32099 |
| n4-Basic | 3.2kb | promoter+enhancer | pGLBasic | BglII/MluI | 35319-32099 |
| n2r-Enhancer | 570bp | promoterreversed | pGLEnhancer | BglII/KpnI | 36392-36964 |
| n3r-Enhancer | 619bp | promoterreversed | pGLEnhancer | BglII/MluI | 32099-32729 |
| dk1-Enhancer | 572bp | promoter | pGLEnhancer | BglII/KpnI | 20950-21522 |
| dk3-Enhancer | 570bp | promoter | pGLEnhancer | BglII/MluI | 32232-32802 |
| dk4-Basic | 1.6kb | promoter+enhancer | pGLBasic | BglII/MluI | 31220-32802 |
| n2-Basic | 570bp | promoter | pGLBasic | BglII/KpnI | 36964-36392 |
| n3-Basic | 619bp | promoter | pGLBasic | BglII/MluI | 32729-32099 |
| n2r-Basic | 570bp | promoterreversed | pGLBasic | BglII/KpnI | 36392-36964 |
| n3r-Basic | 619bp | promoter | pGLBasic | BglII/MluI | 32099-32729 |
| dk1-Basic | 572bp | promoter | pGLBasic | BglII/KpnI | 20950-21522 |
| dk3-Basic | 570bp | promoter | pGLBasic | BglII/MluI | 32232-32802 |
| dk1r-Enhancer | 572bp | promoterreversed | pGLEnhancer | BglII/KpnI | 21522-20950 |
| dk1r-Basic | 572bp | promoterreversed | pGLBasic | BglII/KpnI | 21522-20950 |
| dk3r-Enhancer | 570bp | promoterre versed | pGLEnhancer | BglII/MluI | 32802-32232 |
| dk3r-Basic | 570bp | promoterreversed | pGLBasic | BglII/MluI | 32802-32232 |

Humancelllines293andHepG2weredeterminedtoexpressthegenesofinterest

byanalysisofcellularcDNAwithgenespecificprimers,andw        ereconsequentlyplatedin

96wellformatforthetransienttransfectionluciferaseassays.LNCaP.FGC,whichdid

notshowexpressionofanyofthesegenes,wasusedtotestluciferaseassayresultsina

non-expressingcellline.

Constructsn1andn2



**Figure 9:Constructsn1andn2.**

Constructn2ispredictedbyFirstEFtobeapotentialupstreampromoterforthe
*NPHS1*gene.Thissequenceispositionedunusuallyforapromoterinthatitlieswithin
the *DKFZp564A1164*transcriptio nunit(Figure9).Then1constructincludesthen2
promoterregionplusflanking630bpofupstreamconsensussequencethatwas
consideredapossibleenhancerregion.Then2promoterwasdirectionallysubclonedinto
thepGL -Enhancerand -Basicvectorst otestexpressioninaconstructcontainingand
lackingtheSV40enhancer,respectively.Then1regionwassubclonedintothepGL              -
Basicvectoronly.Then2rpromoterconstructisidenticalton2exceptsubclonedinthe
reverseorientationintopGL   -Enhancerand  -Basicvectors;thisconstructwasdesignedas
apossiblecontrolforn2.IntactpGL       -Enhancerand  -Basicvectors,whichlackedany
insert,wereusedasnegativecontrols,andthepGL        -Controlvectorwhichcontainsan
SV40promoterandenhancerwasu      sedasanexampleofastrongpositive.

Surprisingly, the results of the transfection assay indicated that both n2 and n2r have strong promoter activity in the pGL -Enhancer vector transfected into the HepG2 cell line. The same constructs also show promoter activity in the 293 cell line similar to the positive control, and no activity in LNCap (Figure 10). The strong positive n2 sequence indicates that it may be an alternative upstream promoter for the NPHS1 gene as predicted by FirstEF, and the fact that n2r acts as a strong promoter in both cell lines indicates that it is a bi-directional promoter. These data suggest that the n2r sequence may also function as a downstream alternative promoter for the *DKFZp564A1164* gene. Then1 -Basic construct displayed a reduction of promoter activity compared to n2 -Basic suggesting there may be a silencer in this region causing repression of expression in the cell lines used for this study. Often silencers causing repression of expression are found in the 5′ upstream region of genes (Kemp *et al.*, 2002; Kraner *et al.*, 1992).



**Figure 10: Fold change in relative light units (RLU) of n1 and n2 constructs transfected into 293, HepG2 and LNCap cell lines.**

Constructsn3andn4



**Figure 11:Constructsn3andn4.**

Constructn3isalsopredictedbyFirstEFtobeanalternativepotentialupstream

promoterforthe *NPHS1*gene.Agrowingbodyofdatasuggeststhatmanygenesuse

alternatepro motersindifferenttissues(Asnagli *etal.* ,2002) .Then4constructincludes

then3promoterregionplusflanking2581bpofupstreamconsensussequencethatwas

consideredapossibleenhancerregionforthispromoter(Figure11).Then3 promoter

wasdirectionallysubclonedintothepGL -Enhancerand- Basicvectorstotestexpression

inaconstructcontainingandlackingtheSV40enhancer.Then4regionwassubcloned

intothepGL -Basicvectoronly.Then3rpromoterisidenticalton3exc eptsubclonedin

thereverseorientationintopGL -Enhancerand- Basicvectors.IntactpGL -Enhancerand

-Basicvectors,whichlackedanyinsert,wereusedasnegativecontrols,andthepGL -

ControlvectorwhichcontainsanSV40promoterandenhancerwasuse dasanexample

ofastrongpositive.

Then3promoterregionwaspredictedbyFirstEFtobeapotentialbi　　-directional promoterfor *NPHS1*and *DKFZp564A1164*anddoesshowhigherlevelsofexpressionin thereverseorientation(n3r　-Basicandn3r　-Enhancer)wh encomparedtothenegative controls.However,theresultsofthetransfectionassayindicatethatn3risaweak promoterincomparisontothepositivecontrol(Figure12).Whenthescaleisdecreased inthegraphsothatdifferencesinpromoteractivit　　ycanbevisualizedforthetestregions, a4(n3r　-Basic)and6(n3r　-Enhancer)foldincreaseinexpressionisclearlyvisible(Figure 12and13).Itshouldbenotedthatthepositivecontrolusedintheseexperimentswas suppliedbyPromegaandcontainsa　　verystrongSV40promoterandenhancer,andmost humanpromoterswillnotbeasstrongorstrongerthanthepositivecontrol.Expression oftheforward *NPHS1*constructs,n3　-Basicand　–Enhancerwerebarely1foldgreater thanthenegativecontrols,indica　tingthatn3isprobablynotapromoterforthe　　　*NPHS1* gene.Then4　-Basicconstructreducedpromoteractivitytothatseeninthenegative controlssuggestingtheremaybeasilencerinthisregioncompletelyshuttingoff expression.

Thedifferenceinexpr　essionbetweencelllinesshouldbenotedaswell.While thepreviousconstructsalwayshadhigherexpressionintheHepG2cellline,then3r promotershowsdeferentialexpressiondependingonthevector.Expressionn3rinthe Basicvectorwashigherin　　the293celllinewhileexpressionintheEnhancervectorwas higherintheHepG2cellline.Thismayjustbeanartifactofthelowexpressionlevels, oraninstanceofenhancercompetition.AstudybyG.I.R.Adam　　　*etal.*showedthatthe SV40enhancerus　edinmanyplasmidsfortransienttransfectionassayscanbeastrong

competitorforpositiveandnegativeregulatoryfactorsinacell        -type-specificmanner

(Adam*etal.*,1996) .AlthoughtheSV40enhancerclearlyperformswellinmostof        the

celltypeswehaveexamined,thisfactormayhelpexplainthedifferencesinluciferase

levelsseeninsomecells.



**Figure 12:Foldchangeinrelativelightunits(RLU)ofn3andn4constructstransfe        ctedinto293, HepG2andLNCapcelllines.**

**Figure 13:Foldchangeinrelativelightunits(RLU)ofn3andn4constructstransfectedinto293, HepG2andLNCapcelllinesincomparisontonegativecon        trols(positivecontrolremoved).**

Constructsdk1,dk3anddk4



**Figure 14:Constructsdk1,dk3anddk4.**

Thedk1anddk3promoterswerepredictedbyFirstEFtobepotentialupstream promotersfor *DKFZp564A1164*(Figure14).Th edk4constructincludesthedk3 promoterregionplus1kbofflankingupstreamconsensussequencethatwasconsidereda possibleenhancerregion.Theenhancerregion(dk2)flankingdk1wasnotsubcloneddue todifficultiesinPCRofthisGC -richregion.

Forwardpromoterconstructsdk1 -Enhanceranddk3 -Basicexpressluciferaseatmore than25timesthatofthenegativecontrolsintheHepG2cellline(Figure15).Inthe293 celllinedk1 -Basicexpressesthehighestlevelofluciferaseatalmost10 -foldre lative lightunits(RLU).Againthepromoter/enhancerconstruct,dk4 -Basic,showsareduction inluciferaseactivitycomparedwiththepromoteronlyconstructs,suggestingasilencer maybepresent.

**Figure 15:Foldchangeinrelativelightunits(RLU)ofdk1,dk3anddk4constructstransfectedinto 293,HepG2andLNCapcelllinesincomparisontocontrols.**

Constructsdk3vs.n3

Thedk3andn3predictedpromotersoverlapbyapproximately470bpandeach

extendsbeyondthiscoreregionbyabout100bp.Theorientationofn3anddk3are

oppositetoeachother,whereasn3isinthesameorientationasdk3r,anddk3isinthe

sameorientationasn3r(Figure16).

Expressionlevelswerehighestforthedk3co      nstructs(5to15  -foldincreases).In

spiteoftheoverlapregion,then3rconstructsonlyshoweda4          -foldincreaseinexpression

(Figure17).Thedk3randn3constructshadthelowestexpressionlevels,similartothe

negativecontrols.

**Figure 16: Overlap region of the n3 and dk3 promoters.**



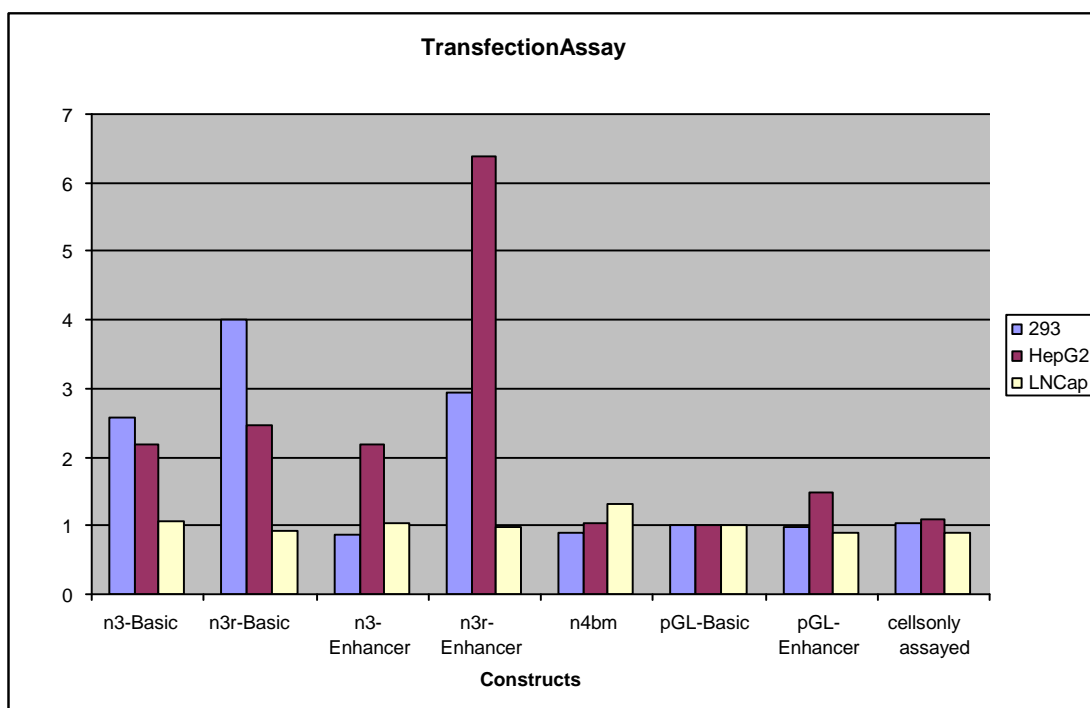**Figure 17: Fold change in relative light units (RLU) of dk3 and n3 constructs transfected into 293 and HepG2 cell lines in comparison to    negative controls.**

Inordertoclarifytheexpressionpatternsseeninthisregion,largerandsmaller constructsweredesignedinthe470bpoverlapregion(Figure18).Intheluciferaseassay thehighestlevelofactivitywasseeninthenewlargerdk      3-Basicconstructindicatingthe extraregioncontainsapowerfulenhancerdrivingthispromoter(Figure19).Thelarger dk3-Enhancerconstructdidnot,however,expressluciferaseathigherlevelsthanthe originaldk3 -Enhancer.SincethepGL    -Enhancer vectorcontainsanSV40enhancer,it maybecompetingforregulatoryfactors,preventingthemfrombindingtotheinsert DNA (Adam*etal*,1996)    .

Thesmallerconstructsandthelargern3constructs,aswellastheoriginaln3,n3r anddk3r constructsallshowedlowlevelsofluciferaseexpressionsimilartothenegative controls.Thesedataseemtoindicatethattheworkingpromoteriswithindk3forward constructandonlyoperatesinonedirectioni.e.isnotbi      -directionalaspredictedb   y FirstEF.Additionally,atleasttwostrongenhancersarelocated5      ′ofthispromoteras evidencedbythehighluciferaseexpressioninthedk3      -Basicandlargedk3  -Basic constructs.Thestrongputativeliverandkidneyenhancersinthisregiondeservef          urther studyincludingthepossibilitythattheSV40enhancermaybecompetingforthesame transcriptionfactors.

**Figure 18: New large and small constructs, their orientation and average RLU.**

**Figure 19: The bight blue and pink bars show the expression patterns of the new large and small constructs.**

## rVISTA Analysis

rVISTA analysis, which detects transcription factor binding sites (TFBS) by clustering and analysis of conserved interspecies sequence, was performed on construct sequences from n1, n2, n3, n4 and dk3 and dk4 using a core similarity of 0.85 and matrix similarity of 0.9, slightly higher than the default parameters. Construct dk1's similarity standards were left at default, 0.75 and 0.8, respectively. All conserved or aligned TFBS that were found in the sequences are listed in table 6A in the appendix.

The results found 11 conserved TFBS in n1 (enhancer region only) including 1 GATA site and 5 CAP sites. All 11 TFBS are found within a 35 bp region immediately 5′ of the promoter. Promoter n2 contained 12 aligned TFBS including 8 CAP sites and 1 each of CETS1P54, ZIC3, CDXA and MZF1.

Promoter n3 contained 9 conserved TFBS located at the 3′ end, and n4 (promoter + enhancer construct) contained 10 conserved TFBS found in clusters throughout the enhancer region. It has been shown that when multiple cis DNA elements are clustered in a region they may work cooperatively to regulate expression (Belsham & Mellon, 2000; Liu *et al.*, 2003). Some of the transcription factor binding sites found in this region included 22 CAP sites, 7 STAT sites, 8 PAX2 sites, 2 GATA sites and 2 YY1 sites.

In the dk1 promoter, only 2 conserved TFBS were found, CAP and ZP1, and both were located in the 3′ end of the promoter. In the dk3 promoter, 4 transcription factor binding sites were found at the 5′ end including 2 PAX2 sites. Recalling that dk3 and n3 overlap by 470 bp, they also share 4 TFBS, and an additional 5 sites are found in the n3 region of this promoter. When the larger promoter construct incorporating all of n3 and dk3 was assayed the results showed very strong expression in the dk3 orientation in the pGL-Basic vector only, suggesting that the extra TFBS found in the n3 region may actually be enhancers for the dk3 promoter. Twenty TFBS were found throughout in the enhancer region of dk4 including 2 PAX2 sites, 4 STAT, 2 GATA and 6 CAP sites. The high concentration of conserved TFBS, especially the clustering of multiple copies of some TFB sites, are consistent with the predicted enhancer role for this *DKFZp564A1164* region.

**DISCUSSION**

Theresultsconfirmthatcomparativesequenceanalysisbetweendivergentspecies

suchhumanandmouseisapowerfultoolforident ifyingregulatoryelementsinnon -

codingconservedsequence (Loots*etal*,2000) .Inthisstudyweusedthewealthof

conservedsequencedataforHSA19andmousetolocateputativepromoterelements,and

exploredtheuseofcomparativese quenceanalysisprogramssuchasPipMakerorVISTA

andthecomputationalpromoterfindingprogram,FirstEF,toassistinlocatingpotential

promoters.ThisisthefirststudydesignedtotesttheFirstEFpredictions,andtheresults

showthat3outof4 predictedpromoterswerefunctionalintheluciferaseassay(Figure

20).However,muchlargernumbersofFirstEFpredictionsneedtobeassayedtoassess

thismethod.



**Figure 20:ThreeoffourFirstEFpredictedpromotersshowede xpressionintheluciferaseassay.
Onepredictedpromoter(n2/n2r)wasfoundtohaveexpressioninbothorientations,althoughitwas
notpredictedbyFirstEFtobebi -directional.**

Theresultsalsodemonstratethattestingpotentialregulatoryelementsi n

transientlyexpressedluciferasereporterconstructstransfectedintoculturedmammalian

celllinesisareliablemethod,andbecomesahighthroughputmethodwhenperformedin

a96wellformat.Inthisstudythechoiceofcelllinewasfoundtobeof                critical

importancetoassayresults.ForexampleingeneraltheHepG2celllineproducedhigher

luciferasevalueswhentransfectedwiththeseparticularpromoters.However,in3

constructs293showedhighervalues:n3Basic,n3rBasicanddk4Basic.The                LNCapcell

line,ontheotherhand,wasapoorreporteralltogether.Manygenesusealternativestart

sitesandpromotersindifferenttissues,sopromotersshouldbetestedinatleasttwo

differentcelllinesthatarebasedontheresultsofprelimina        ryexpressiondata  (Asnagli*et
al.*,2002) .

Theresultsdidnot,however,showthat        *NPHS1,HSPOX1*  and*DKFZp564A1164*

shareasinglebi   -directionalpromoter(n3/dk3).Then3constructdoesnothavepromoter

activityinthecelllineswetest        edand,therefore,isprobablynotapromoterfor        *NPHS1*

and *HSPOX1*.However,n2unexpectedlyturnedouttobeabi        -directionalpromoter.

Then2constructisanexcellentexampleoftheimportanceoftestingall

hypotheticalpromotersinbothorientations.    Although,n2waspredictedbyFirstEFtobe

apromoteronlyinonedirectionforthe        *NPHS1*geneitexpressedhighlevelsof

luciferaseactivityinbothorientationsindicatingitisastrongbi        -directionalpromoter.

Sequencelengthwasalsoshowntobeo        fimportanceinthisstudybecause

althoughn3anddk3share470bpofsequencewitheachotheritwasthe70to100base

pairsthattheydidnotsharethatwasfoundtoenhanceorreduceexpression.The

rVISTAdatashowedthatseveralpotentialtranscrip        tionfactorbindingsites(TFBS)exist

ontheperipheryofthiscorepromoterregion.Whenalargerconstructwasdesignedit

wasrevealedthatthedk3promoterwasfurtherenhancedwhilethen3directionremained

thesamesuggestingthattheextrabasep        airscontainTFBSthatactasenhancersonthe

dk3promoter.

Throughoutthisprojecttheexperimentswereformattedtoestablishthe

technologyandmethodsforahighthroughputassayofpromoterandenhancerelements.

UsingFirstEFasaguide,putativep        romoterscanbequicklyassayedforactivity.Inour

study3of4promotersassayedshowedconsiderableincreasesinluciferaseactivityover

negativecontrolsdenotingaworkingpromoter(Figure20).Although,constructn3only

showedaslightincrease    inluciferaseactivityoverthenegativecontrols,dk3,which

overlapsthesameregion,butwasclonedintheoppositeorientation,showedverystrong

luciferaseactivity.Notallmammalianpromotersaregoingtobeasstrongaspositive

controlsandwes    houldexpecttoseeahighdegreeofvariabilityinexpression.

Theconsensussequenceupstreamofeachpredictedpromoterwastestedfor

enhanceractivityinthepGL3    -Basicvector,andall3ofthese"enhancer"constructshad

reducedactivityrelativeto    theshorterpromotersequences.Aspreviousstudieshave

shownthe5 ′regionofapromotercancontainsilencersitescausingtranscriptional

repression (Kemp*etal.* ,2002;Kraner *etal.* ,1992) .Whentheseregionswereexamined

byrVISTA  anumberoftranscriptionfactorbindingsiteswerefound,someofwhichare

knowntoberepressorsforcertaingenesorinspecifictissues.ForexampleYY1,PAX2

andCIZbindingsiteswerefoundinoneormoreoftheenhancerregionsandallhave

beens howntoreduceexpressioninpreviousstudies        (Havik*etal.* ,1999;Kim *etal.* ,

2003;Shen *etal.* ,2002) .

TheenhancerconstructsweremoredifficulttoPCRandsubcloneduetotheir largersizeandhighGCcontentmakingthegenerationof theseregionsmoretime consuming.Forthisreasonhighthroughputassaysofputativeenhancer/silencerregions maynotbeabletokeeppacewithassaysofthepromoterregions.

The5 ′SMARTRACEexperimentsperformedtoidentifythepredictedfirstexon ofeach gene,werealsoaslowingpointinthishighthroughputpipeline.Confirmingthefirst exonisgoingtobecriticalforprovingwhichgenethesepromotersoperateon,however, itmaytakemoretime.TryingdifferentRACEkitsoralternativelyamp lifyingRT -PCR productswithgenespecificprimersmightyieldbetterresults.Ofthefourpromoters testedonlyone,dk3,waslocatedadjacenttothefirstexonofagene( *DKFZp564A1164*). Theother3promotersare6to10kbawayfromtheknowntranscri bedsequencesofthe genestheyarepredictedtooperateon.

ThefailuretoRACE*NPHS1*and *HSPOX1*couldalsobetaken,togetherwith reporterresults,toindicatethatFirstEFfailedtofindeithergene'spromoterandthatthe predictionofupstreamexons maybeincorrect.InthiscaseitisclearthatFirstEFdidfail topredictwhatappeartobethemostcommonlyusedfirstexonsforboth *NPHS1* and *HSPOX1*.OneofthepurposesofthisstudywastoprovidedatatotestFirstEF predictionsandfeedbackth eresultstoFirstEF'screators.BecauseFirstEFisarelatively newprogram,suchfeedbackwillbehelpfulinrefiningitspredictionalgorithms.

Then2construct,whichwaspredictedtooperateinoneorientationandtoprovide apotentialupstreamprom oterfor *NPHS1*gene,wasfoundtobeastrongbi -directional promoter.Theclosestgenethatthereverseorientationofn2(n2r)couldbeoperatingon

is *APLP1,* 8Kbaway. However, itismorelikelythatn2rcouldbeaninternalalternative promoterfor *DKFZp564A1164,* andmaypotentiallydefinealternativestartsitesforboth *NPHS1* and *DKFZp564A1164*. Furtherexperimentsarenecessaryinthisregionto confirmwhichgenesthesepromotersareoperatingon.

## Conclusion

Forthisstudy, ahighthroughputmet hodforidentifyingandtestingregulatory elementswasexamined. Inaddition, thevalidityofpromoterspredictedbyFirstEFwas tested. Itwasfoundthatbycombiningcomputerbasedpromoterandfirstexon predictionsfromFirstEF (Davuluri*etal.* ,2001) withPCR -basedcloningtogenerate luciferasereporterconstructs, andbytestingreporteractivityinculturedmammaliancells platedina96wellformatonecouldidentifypromoteractivityinarelativelyhigh throughputmanner.

The datageneratedinthisstudysuggestthatFirstEFpredictionsaresometimes incorrect. Therefore, havingastrategyfordefiningwhichFirstEFpredictedpromotersto testfirstmayacceleratetheprocess. Initiallytestingpromotersthatareataconfirm ed transcriptionstartsiteforagene, atapossiblealternatetranscriptionstartsiteorina regionofconservedsequencewouldbethebestcandidates, whilepromoterspredictedin genedesertregionsmaynotbeaseasytoconfirm.

Theluciferaseassay lentitselfverywelltothehighthroughputsearch, however thesubcloningdidnotalwaysgosmoothly. Thenumerousstepsthatthistraditional subcloningmethodrequiresweretimeconsumingandincreasedtheopportunitiesfor

errors. A faster method that skips many of the traditional subcloning steps, such as the Creator™ system by Clontech is currently being investigated by our lab.

The development and testing of substantially larger enhancer/silencer regulatory elements may not be possible at this time using these high throughput methods. These regulatory elements are generally GC rich making them more difficult to PCR and subclone. Additionally, confirming upstream untranslated first exons was not possible within this timescale using the SMARTR ACE protocol. It will be necessary to further explore the limitations within these procedures in order to confirm these and future regulatory elements. Alterations and modifications to these protocols, as well as investigating new techniques may be neces sary.

# REFERENCES

AdamG.I.,MillerS.J.,UllerasE.,andFranklinG.C.(1996).Cell          -type-specific
    modulationofPDGF  -Bregulatoryelementsviaviralenhancercompetition:a
    caveatfortheuseofreferenceplasmidsintransienttransfe           ctionassays. *Gene*
    **178:**25- 29.

AsnagliH.,AfkarianM.,andMurphyK.M.(2002).IdentificationofanAlternative
    GATA-3PromoterDirectingTissue-   SpecificGeneExpressioninMouseand
    Human. *JImmunol*  **168:**4268 -71.

BelshamD.D.,andMellonP.L.(2000        ).TranscriptionfactorsOct  -1andC/EBPbeta
    (CCAAT/enhancer-bindingprotein -beta)areinvolvedintheglutamate/nitric
    oxide/cyclic-guanosine5' -monophosphate-mediatedrepressionofmediated
    repressionofgonadotropin -releasinghormonegeneexpression.*Mo     l.Endocrinol.*
    **142**12-  228.

BrayN.,Dubchak,I.andPachter,L(2003).AVID:AGlobalAlignmentProgram.
    *GenomeResearch*  **13:**97.

ChakravartiA.(2002).Acompellinggenetichypothesisforacomplexdisease:
    PRODH2/DGCR6variationleadstoschizophrenia     susceptibility. *Proc.Natl.
    Acad.Sci.USA,*  **99:**4755  -4756.

DavuluriR.V.,GrosseI.,andZhangM.Q.(2001).Computationalidentificationof
    promotersandfirstexonsinthehumangenome.        *NatGenet* **29:**412-  7.

DehalP.,PredkiP.,OlsenA.S.,Kobayashi       A.,FoltaP.,LucasS.,LandM.,TerryA.,
    EcaleZhouC.L.,RashS.,ZhangQ.,GordonL.,KimJ.,ElkinC.,PollardM.J.,
    RichardsonP.,RokhsarD.,UberbacherE.,HawkinsT.,BranscombE.,and
    StubbsL.(2001).Humanchromosome19andrelatedregionsin      mouse:
    conservativeandlineage -specificevolution. *Science* **293:**104 -11.

DubchakI.,BrudnoM.,LootsG.G.,PachterL.,MayorC.,RubinE.M.,andFrazerK.
    A.(2000).Activeconservationofnoncodingsequencesrevealedbythree        -way
    speciescomparisons. *GenomeRes*  **10:**1304  -6.

GoodmanB.,RutbergJ.,LinW.,PulverA.,andThomasG.(2000).Hyperprolinaemiain
    patientswithdeletion(22)(q11.2)syndrome.        *JInheritMetabDis*  **23:**847- 848.

HallmanN.,HjeltL.,andAhvenainenE.K.(1956).Nephroticsyndro           meinnewbornand
    younginfants. *AnnPaediatFenn*  **2:**227- 241.

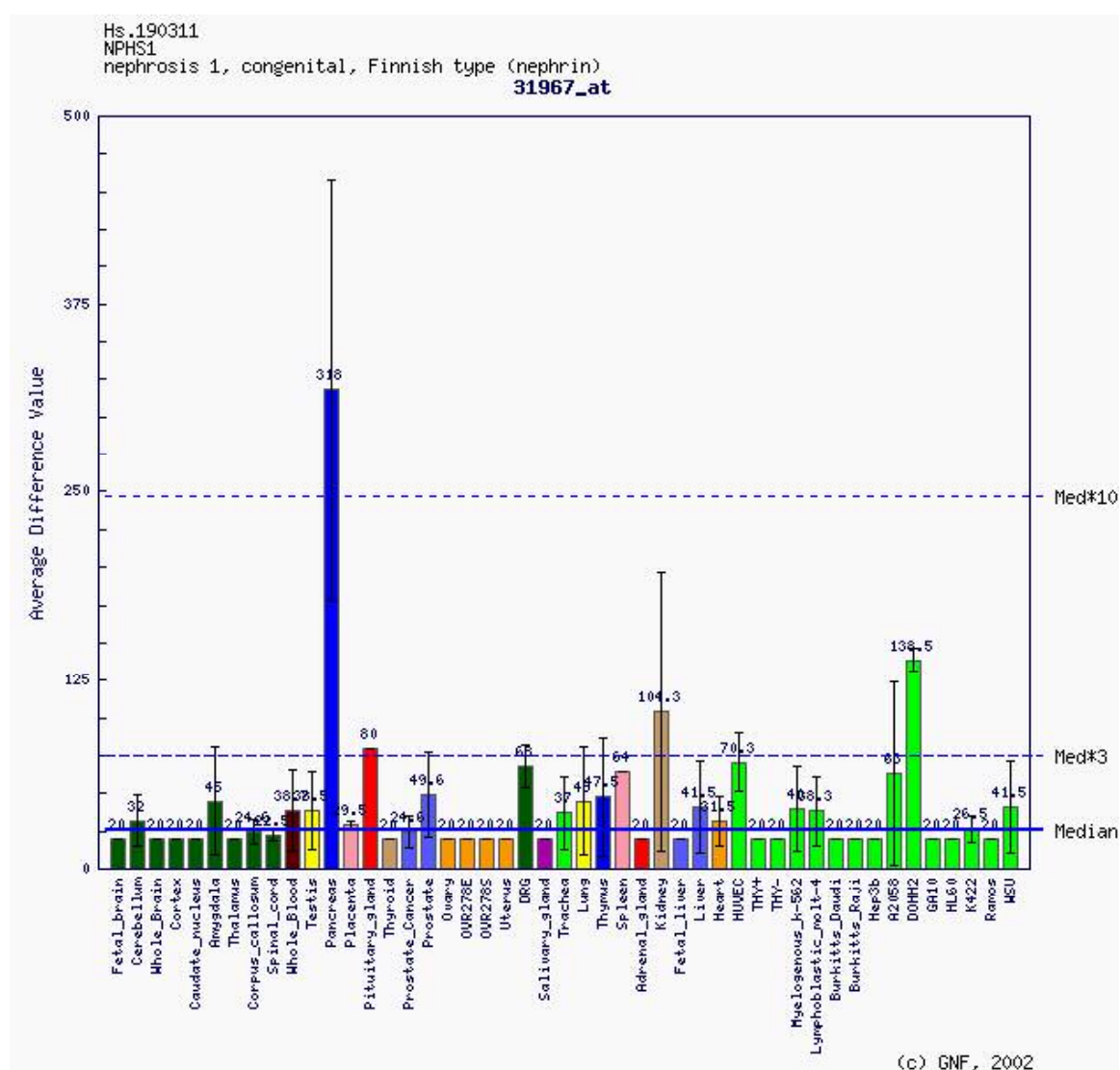HardisonR.C.,OeltjenJ.,andMillerW.(1997).Longhuman    -mousesequence alignmentsrevealnovelregulatoryelements:areasontosequencethemouse genome. *GenomeRes* **7:**959 -966.

HavikB., RagnhildstveitE.,LorensJ.B.,SaelemyrK.,FauskeO.,KnudsenL.K.,and FjoseA.(1999).AnovelpaireddomainDNArecognitionmotifcanmediatePax2 repressionofgenetranscription. *BiochemBiophysResCommun.* **266:**532- 541.

HoodL.,KoopB.F.,Row  enL.,andWangK.(1993).HumanandmouseT    -cell-receptor loci:theimportanceofcomparativelarge  -scaleDNAsequenceanalyses. *Cold SpringHarbSympQuantBiol* **58:**339- 348.

IhalmoP.,PalmenT.,AholaH.,ValtonenE.,andHolthoferH.(2003).Filtrini    sanovel memberofnephrin -likeproteins. *BiochemBiophys* **300:**364 -370.

KempD.M.,LinJ.C.,UbedaM.,andHabenerJ.F.(2002).NRSF/RESTconfers transcriptionalrepressionoftheGPR10geneviaaputativeNRSE/RE    -1located inthe5'promoterregion. *FEBSLett* **531:**193- 8.

KestilaM.,LenkkeriU.,MannikkoM.,LamerdinJ.,McCreadyP.,PutaalaH., RuotsalainenV.,MoritaT.,NissinenM.,HervaR.,KashtanC.E.,PeltonenL., HolmbergC.,OlsenA.,andTryggvasonK.(1998).Positionallyclonedgenefor    a novelglomerularprotein --nephrin--ismutatedincongenitalnephroticsyndrome. *MolCell* **1:**575- 82.

KimJ.,KollhoffA.,BergmannA.,andStubbsL.(2003).Methylation    -sensitivebinding oftranscriptionfactorYY1toaninsulatorsequencewithinthepat    ernally expressedimprintedgene,Peg3. *HumMolGenet.* **12:**233 -245.

KranerS.D.,ChongJ.A.,H.J.T.,andMandelG.(1992).SilencingthetypeIIsodium channelgene:amodelforneural    -specificgeneregulation. *Neuron* **1:**37- 44.

LenkkeriU.,MannikkoM  .,McCreadyP.,LamerdinJ.,GribouvalO.,NiaudetP.M., AntignacC.K.,KashtanC.E.,HombergC.,OlsenA.,KestilaM.,and TryggvasonK.(1999).Structureofthegeneforcongenitalnephroticsyndromeof thefinnishtype(NPHS1)andcharacterizationof    mutations. *AmJHumGenet* **64:** 51-61.

LiuM.,LeibowitzJ.L.,ClarkD.A.,MendicinoM.,Q.N.,DingJ.W.,D'AbreoC.,Fung L.,MarsdenP.A.,andLevyG.A.(2003).Genetranscriptionoffgl2in endothelialcellsiscontrolledbyEts    -1andOct -1andre quiresthepresenceof bothSp1andSp3. *EurJBiochem.* **270:**2274 -2286.

LootsG.G.,LocksleyR.M.,BlankespoorC.M.,WangZ.E.,MillerW.,RubinE.M.,
        andFrazerK.A.(2000).Identificationofacoordinateregulatorofinterleukins4,
        13,and5byc ross-speciessequencecomparisons. *Science* **288:**136 -40.

LootsG.G.,Ovcharenko,I.,Pachter,L.,Dubchak,I.,Rubin,E.(2002).rVISTAfor
        comparativesequence -baseddiscoveryoffunctionaltranscriptionfactorbinding
        sites. *GenomeRes.* **12**832 -9.

MatherJ.P.,andRobertsP.E.(1998)."Introductiontocellandtissueculture:theoryand
        technique,"PlenumPress,NewYork.

MayorC.,BrudnoM.,SchwartzJ.R.,PoliakovA.,RubinE.M.,FrazerK.A.,PachterL.
        S.,andDubchakI.(2000).VISTA:visuali         zingglobalDNAsequencealignments
        ofarbitrarylength. *Bioinformatics* **16:**1046- 7.

MoellerM.J.,KovariI.A.,andHolzmanL.B.(2000).Evaluationofanewtoolfor
        exploringpodocytebiology:mouseNphs15'flankingregiondrivesLacZ
        expressioninpo docytes. *JAmSocNephrol* **11:**2306 -14.

ParsonsS.J.,RhodesS.A.,ConnerH.E.,ReesS.,BrownJ.,GilesH.(2000).Useofa
        dualfireflyandRenillaluciferasereportergeneassaytosimultaneouslydetermine
        drugselectivityathumancorticotrophinrele         asinghormone1and2receptors.
        *AnalyticalBiochemistry* **281:**187- 192.

PrazV.,PerierR.,BonnardC.,andBucherP.(2002).TheEukaryoticPromoter
        Database,EPD:newentrytypesandlinkstogeneexpressiondata.         *NucleicAcids
        Res* **30:**322 -324.

SchwartzS.,ZhangZ.,FrazerK.A.,SmitA.,RiemerC.,BouckJ.,GibbsR.,Hardison
        R.,andMillerW.(2000).PipMaker      --awebserverforaligningtwogenomicDNA
        sequences. *GenomeRes* **10**577- 86.

ShenZ.J.,NakamotoT.,TsujiK.,NifujiA.,MiyazonoK.,Komori         T.,HiraiH.,and
        NodaM.(2002).Negativeregulationofbonemorphogeneticprotein/Smad
        signalingbyCas -interactingzincfingerproteininosteoblasts.         *JBiolChem.* **277:**
        29840-6.

SherfB.A.,NavarroS.L.,HannahR.R.,WoodK.V.(1996).Dual          luciferasereporter
        assay:anadvancedco -reportertechnologyintegratingfireflyandrenilla
        luciferaseassays. *In*"PromegaNotesMagazine",pp.02.

WiemannS.,WeilB.,WellenreutherR.,GassenhuberJ.,GlasslS.,AnsorgeW.,Bocher
        M.,BlockerH.,Bauersachs S.,BlumH.,LauberJ.,DusterhoftA.,BeyerA.,
        KohrerK.,StrackN.,MewesH.W.,OttenwalderB.,ObermaierB.,TampeJ.,

Heubner D., Wambutt R., Korn B., Klein M., and Poustka A. (2001). Toward a catalog of human genes and proteins: sequencing and analy sis of 500 novel complete protein coding human cDNAs. *Genome Res* **11:**422- 35.

Zhou Z. -Q., and Walter C.A. (1998). Cloning and characterization of the promoter of baboon XRCC1, a gene involved in DNA strand -break repair. *Somat Cell Mol Genet.* **24:**23- 39.

**APPENDIX**

**Figure1A:Microarrayexpressiondata**



Microarrayexpressiondatafrom GeneExpressionAtlas( http://expression.gnf.org/cgi-bin/index.cgi)microarraydatabasefor *NPHS1*showinghi ghestexpressioninthe pancreas.

**APPENDIX**

**Figure2A:Microarrayexpressiondata**



Microarrayexpressiondatafrom GeneExpressionAtlas( http://expression.gnf.org/cgi-bin/index.cgi)microarraydat abasefor *HSPOX1*showinghighestexpressioninthe kidney.

**APPENDIX**

**Table 1A:PrimersforcDNAanalysis.**

| Genename | Forwardprimer(exon) | Reverseprimer(exon) |
|---|---|---|
| *NPHS1* | GAGGACCGAGTC AGGAACGAA(26) | CTGCACTTCATCGTA GAGGGGT(28) |
| *DKFZp564A1164* | AGCAAAAGAACC TGATGCGAATC(13) | TTGATGTAGCTG GTGAAAGCTCG(15) |
| *HSPOX1* | CCATGAGGAARCTGT TCGCC(9) | TGCTAGTGGGGT ATCCTTC(11) |
| β−actin | GCGGGAAATCGTGCG TGACATT | GATGGAGTTGAA GGTAGTTTCGTG |

**Table 2A:Tissuearrayprobes.**

| Gene | *HSPOX1* |
|---|---|
| accession number | NM_021232,mRNA |
| forward | GGGCAGTTGGTGAACTTGCT |
| reverse compliment | TCAGCTCTCCTGTGCCCTTA |
| reversew/ **t7** | **TAATACGACTCACTATAGGG**TCA GCTCTCCTGTGCCCTTA |
| | |
| gene | *NPHS1* |
| accession number | NM_004646 |
| forward | GAGGAGGTGTCTTATTCCCG |
| reverse compliment | TCCAGAGTGTCCAAGTCTCC |
| reversew/ **t7** | **TAATACGACTCACTATAGGG**TCC AGAGTGTCCAAGTCTCC |
| | |
| gene | *DKFZP564A1164* |
| accession number | NM_032123 |
| forward | ACTACAAGGTCCGAGGAGTC |
| reverse compliment | TGCCCTGGCTCTGTAAAGTC |
| reversew/ **t7** | **TAATACGACTCACTATAGGG**TGC CCTGGCTCTGTAAAGTC |

**APPENDIX**

**Table 3A:Tissuehybridizationresults.**

|  | *NPHS1* | *HSPOX1* | *DKFZp564A1164* |
|---|---|---|---|
| Lung | - | - | +/ - |
| Skin | + | +/ - | +exterior |
| Muscle,skeletal | - | - | - |
| Heart, muscle | - | - | +/ - |
| Stomach | - | ++ | ++ |
| Esophagus | +/- | +/ - | +/ - |
| Smallintestine | + | +/ - | + |
| Colon | - | - | +/ - |
| Liver | - | +/ - | + |
| Spleen | - | - | - |
| Pancreas | + | - | - |
| Salivarygland | - | + | +/ - |
| Pituitarygland | - | - | - |
| Adrenalgland | - | - | - |
| Thyroidgland | - | - | -(+supportti ssue) |
| Parathyroidgland | - | - | - |
| Thymusgland | +/- | +/ - | ++ |
| Tonsil | + | + | + |
| Bonemarrow | - | - | - |
| Breast | -( +ingland) | -( +ingland) | - ( +ingland) |
| Uterus | ++ | + | ++ |
| Cervix | +/- | - | +/-endothelial |
| Ovary | ++ | ++ | ++ |
| Kidney | +(tubulesonly) | +/ - | +(tubulesonly) |
| Prostategland | ++ | + | ++ |
| Testis | ++ | ++ | ++ |
| Omentum | - | +/ - | + |
| Peripheralnerve | + | +/ - | + |
| Cerebralcortex | ++ | ++ | + |
| Cerebellum | ++ | + | +perkingi |

Table3A:  –negative;+/ -weakpositive;+positive;++strongpositive

**APPENDIX**

**Figure3A:Livertissueslide**



Livertissue: *DKFZp564A1164*t7mRNAprobelabeledwithdigandhybridizedtoa
MaxArraynormalhumantissueslide(ZymedLaboratories,Inc.)redcolor
indicatespositivehybridization

**APPENDIX**

**Figure4A:Livertissueslide**



Livertissue:  *HSPOX1*t7mRNAprobelabeled  withdigandhybridizedtoaMaxArray
     normalhumantissueslide(ZymedLaboratories,Inc.)redcolorindicatespositive
     hybridization

**APPENDIX**

**Figure5A:Livertissueslide**



Livertissue:  *NPHS1*t7mRNAprobelabeledwithdigandhybridizedtoaMaxArray
        normalhumantissueslide(ZymedLaboratories,Inc.)redcolorindicatespositive
        hybridization.
.

**APPENDIX**

**Table 4A:PrimersforPCRofpromoterorpromoter+enhancerconstructs.**

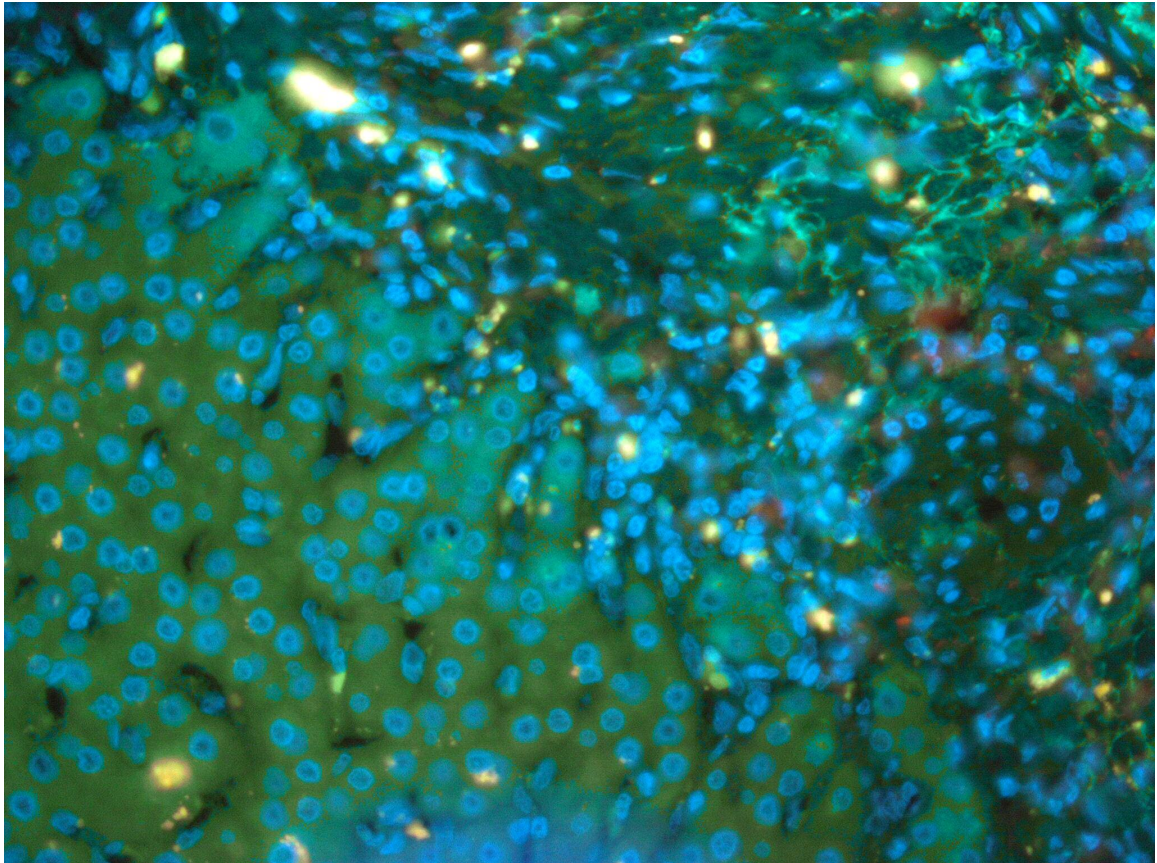| Name | Forward(lowercaselettersare restrictionenzymesequence) | Reverse(lowercaselettersare restrictionenzymesequence) |
|---|---|---|
| n1 | ggaagatctCTGCAGGCA AAGCCGGAGCC | cggggtaccccgAGGTTT GGAGGTCTC |
| n2 | ggaagatctCTGCAGGCA AAGCCGGAGCC | cggggtaccccAAAGGCT GTAACAAAGCC |
| n3 | ggaagatcttccACTCTCTCC CTTCCCTCC | cgacgcgtcgTTCTCGCT AGTGAAGAGGCA |
| n4 | ggaagatcttccACTCTCTCC CTTCCCTCC | cgacgcgtcgTCTCGAAC TCCTGATCTTAG |
| n2r | cggggtaccccTGCAGGC AAAGCCGGAGCC | ggaagatctAAAGGCTGT AACAAAGCC |
| n3r | cgacgcgtcGtcTTCCACTCT CTCCCT TCC | ggaagatctTCTCGCTAG TGAAGAGGCA |
| dk1 | cggggtaccccgAAGGAC GCTCCTGGCGGC | ggaagatcttccAAGGCT GGACAGCTCAGC |
| dk2 | cggggtaccccgTGTGAG AGGGCCCCAGGT | ggaagatcttccAAGGCT GGACAGCTCAGC |
| dk3 | cgacgcgtcgaATTGAGC TGGGGGCGCCCA | ggaagatcttccGGGGCA GCAGGGCTGAGC |
| dk4 | cgacgcgtcgaAATCCTC CTGGGCCTGTG | ggaagatcttccGGGGCA GCAGGGCTGAGC |
| dk1r | ggaagatcttccAAGGACGCT CCTGGCGGC | cggggtaccccgAAGGCT GGACAGCTCAGC |
| dk3r | ggaagatcttccTTGAGCTGG GGGCGCCCA | cgacgcgtcgaGGGGCAG CAGCGGCTGAGC |
| dk3 large | ggaagatcttccACTCTCTCC CTTCCCTCC | cgacgcgtcgaGGGGCAG CAGCGGCTGAGC |
| n3large | ggaagatcttccGGGGCAGCA GGGCTGAGC | cgacgcgtcGtcTTCCAC TCTCTCCCTTCC |
| dk3 small | ggaagatcttccTTGAGC TGG GGGCGCCCA | cgacgcgtcgTTCTCGCT AGTGAAGAGGCA |

**APPENDIX**

**Table 5A:Primersfor5'SMARTRACEof**   *NPHS1*and  *HSPOX1*

*NPHS1*Raceprimers

| Name | Sequence | Size |
|------|----------|------|
| rn1 | GGATGGAGAGGATCACTCTGGGAGACACGA | 30bp |
| rn2 | CCTGAAAACCTGACGGTGGTGGAGGGGGCC | 30bp |
| rn3 | CGGAGTATGAGTGCCAGGTCGGCCGCTCTG | 30bp |

*HSPOX1*RACEprimers

| Name | Sequence | Size |
|------|----------|------|
| rh1 | GGGAACAGAGCACGTAACAGGTCCGGAGC | 29bp |
| rh2 | CTCACCAGCCACAAACTGCCCATAGACGG | 29bp |
| rh3 | ATAGCACCGAGGTTCCCCTCATACCACGCC | 30bp |

**APPENDIX**

**Table 6A:Transcriptionfactorbindingsites(TFBS)foundbyrVISTA**

| Promoter/Enhancer | TFBS | NumberofHits |
|---|---|---|
| n1=11conservedTFBS enhancerregion | AP2ALPHA | 2 |
| | CAP | 5 |
| | GATA | 1 |
| | TEF1_Q6 | 1 |
| | GEN_INI_B | 1 |
| | HOXA4_Q2 | 1 |
| n2=12alignedTFBS promoterregion | CAP | 8 |
| | CETS1P54 | 1 |
| | ZIC3 | 1 |
| | CDXA | 1 |
| | MZF1 | 1 |
| n3=9conservedTFBS promoterregion | CAP | 2 |
| | STAT | 2 |
| | CETS1P54 | 1 |
| | PAX2 | 2 |
| | MYB_Q6 | 1 |
| | SRY | 1 |
| n4=60conservedTFBS enhancerregion | MYB_Q6 | 1 |
| | CAP | 20 |
| | CDXA | 1 |
| | STAT | 5 |
| | PAX2 | 6 |
| | PAX4 | 1 |
| | HOXA4_Q2 | 2 |
| | TEF1_Q2 | 1 |
| | GEN_INI_B | 4 |
| | GATA | 2 |
| | CEBP | 1 |
| | TCF4_Q5 | 1 |
| | CETS1P54 | 2 |
| | NFAT_Q6 | 1 |
| | YY1 | 2 |
| | PEA3_Q6 | 1 |
| | AP2ALPHA | 1 |
| | SPZ1 | 1 |
| | DBP_Q6 | 1 |
| | EN1 | 1 |
| | GR_Q6 | 1 |
| | PU1_Q6 | 1 |
| | NKX62_Q2 | 1 |
| | OCT1 | 1 |
| | CIZ | 1 |

**APPENDIX**

**Table6A** : **Transcriptionfactorbindingsites(TFBS)foundbyrVISTA**

| Promoter/Enhancer | TFBS | NumberofHits |
|---|---|---|
| dk1=2conservedTFBS promoterregion | CAP | 1 |
| | ZP1 | 1 |
| dk3=4conservedTFBS promoterregion | STAT | 1 |
| | PAX2 | 2 |
| | CAP | 1 |
| | MYB_Q6 | 1 |
| | SRY | 1 |
| dk4=20conservedTFBS enhancerregion | PAX2 | 2 |
| | CIZ | 1 |
| | STAT | 4 |
| | LPOLYA_B | 1 |
| | CDXA | 1 |
| | GATA | 2 |
| | CAP | 6 |
| | HSF1 | 1 |
| | AP2ALPHA | 1 |
| | CETS1P54 | 1 |